



# **Robust Speaker Segmentation for Meetings: The ICSI Spring 2006 Diarization System**

Xavier Anguera  
Chuck Wooters  
José Pardo

**RT-06S Meeting Recognition Workshop**

**May 3<sup>rd</sup>, Bethesda, Maryland**



# Outline

- Tasks we participated in for RT-06s
- System description
- What's new since RT-05s
- Eval results
- Post-evaluation analysis
- Future work

# Tasks and Submissions

(23 systems submitted)

## ■ SPKR

### □ Conf room:

- MDM (1p, 4c)
- SDM (1p, 1c)

### □ Lecture room:

- ADM (1p, 3c)
- MDM (1p, 3c)
- SDM (1p, 1c)
- MSLA (1p)

## ■ SAD

### □ Conf room:

- MDM (1p)
- SDM (1p)

### □ Lecture room:

- ADM (1p)
- MDM (1p)
- SDM (1p)

# The ICSI Speaker Diarization System

- Overall goal: Robustness and Portability across domains
- Agglomerative clustering
- Cluster merging uses a modified BIC
  - No penalty term in the BIC formula
- No pre-trained acoustic models
- Take advantage of multiple mics, when they are available, using delay&sum.

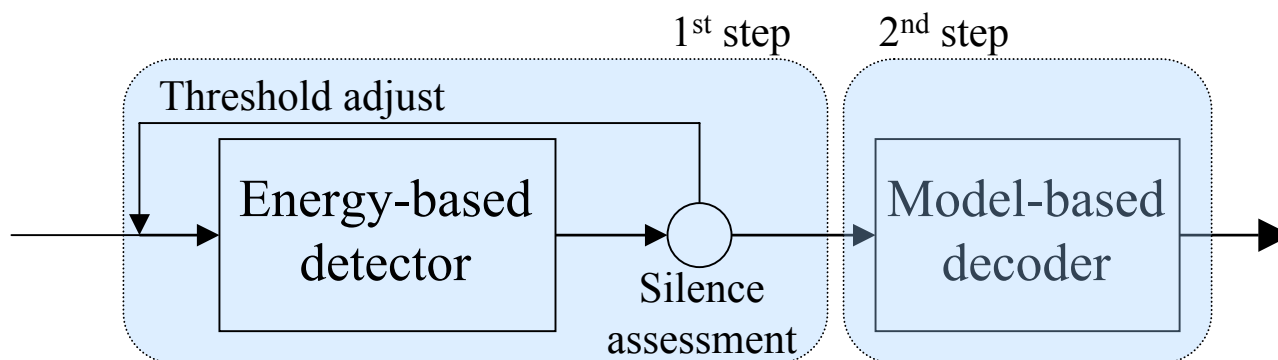
# What's new?

- A new train-free speech/non-speech detector
- Improved cluster initialization (“friends and enemies” algorithm).
- Frame-based purification during merging.
- Initial number of clusters determined semi-automatically.
- Average speaker turn length only acoustically-driven.
- For MDM:
  - Modified delay-sum algorithm.
  - using delays in combination with acoustic features.

# New Speech/non-speech

- Two-pass algorithm:

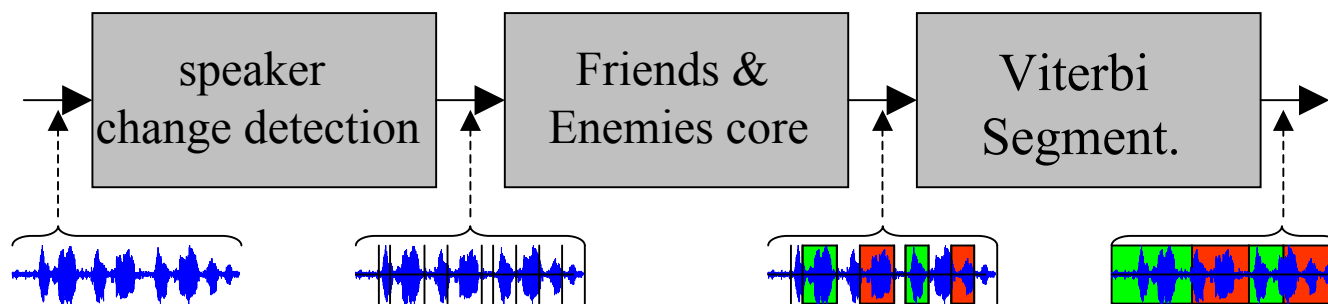
- Step 1: Energy-based detector to locate silence regions with high confidence.
- Step 2: Iterative clustering approach.
  - Models initialized with silence regions from first pass.
  - Convergence based on likelihood.



# Improved Initialization

“Friends and enemies” algorithm:

1. Speaker change detection using BIC.
2. Create a first cluster as the most likely segment according to a general model.
3. Find the friends for this segment with highest cross-likelihood.
4. Iteratively find new “enemies” with smallest cross-likelihood to existing clusters, and find its friends.
5. Train models and segment all the data.



# Frame-based Purification

- Simpler than last year's segment-based purification.
- Problem: BIC-based cluster merging is adversely affected by “impure” data.
  - I.e. non-speech frames that occur in all clusters difficult to discriminate between clusters.
- Observation: most non-speech frames obtain the best lkl in all models.
- Proposed solution: Try to detect “impure” frames, using a likelihood-based metric, and exclude these during cluster comparison.

$$\bar{\mathcal{L}}(x[i] | \Theta_A) = \frac{1}{Q} \sum_{j=-(Q/2)}^{(Q/2)-1} \sum_{m=1}^{\tilde{M}} \log(W_A[m] \mathcal{N}_{A,m}(x[i+j]))$$



# Initial Number of Clusters

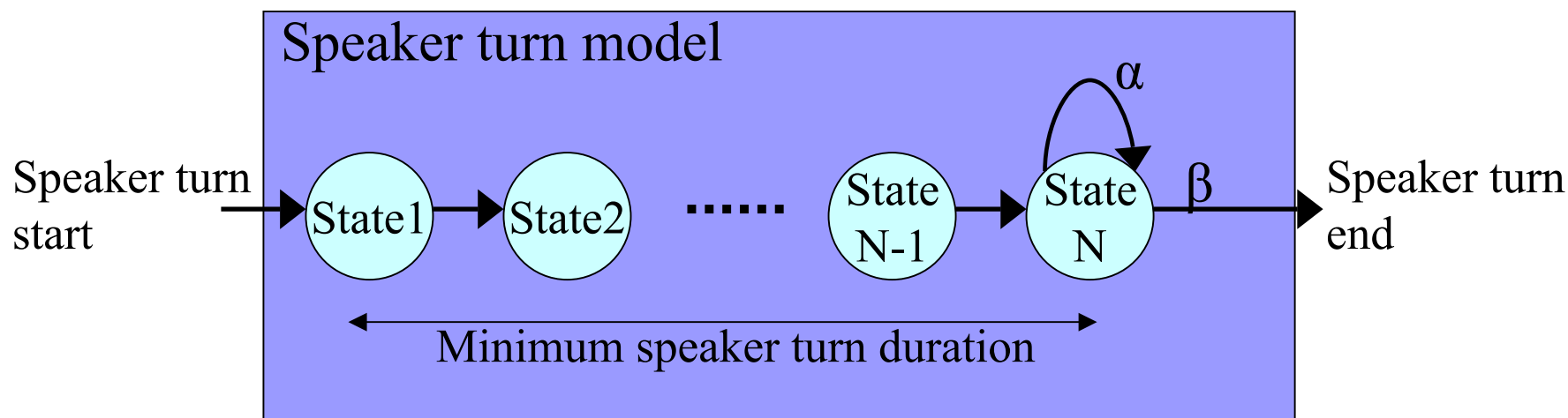
- Last year an informed guess was made about the number of initial clusters (10 for Conf., 5 for Lect.).
- This year the initial number of clusters ( $K_{init}$ ) is determined based on the length of the meeting ( $N_{total}$ ), using:

$$K_{init} = \frac{N_{total}}{GM_{clus} CCR}$$

$GM_{clus}$  = # Gauss per cluster,  $CCR$  = # frames necessary to train a Gaussian, to be optimized.

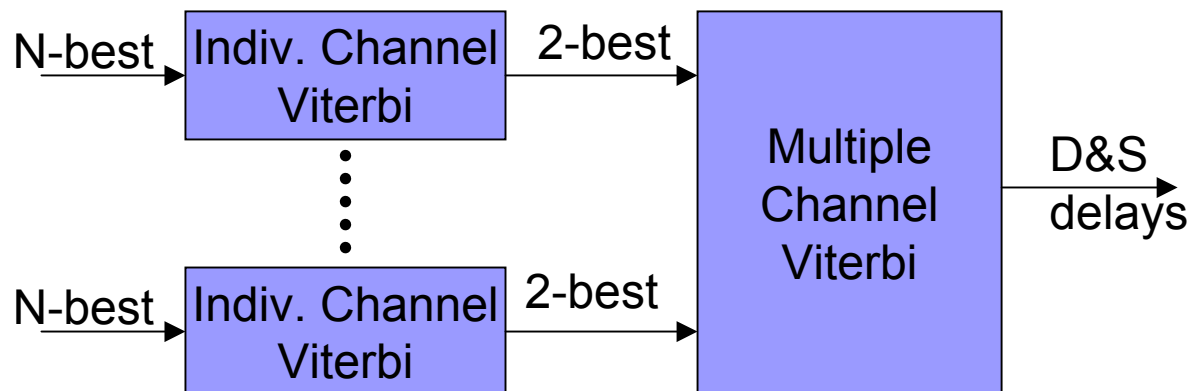
# Modified average speaker turn length

- The average speaker turn duration was artificially constrained by the chosen minimum turn duration and parameters  $\alpha$  and  $\beta$ .
- Making  $\alpha = 1$  and  $\beta = 1$  makes the average turn duration dependent only on the acoustic data.



# MDM-specific Improvements

- Modified delay-sum algorithm (new code)
  - New noise thresholding to eliminate “bad quality” delays.
  - Smoothing of delays using Viterbi in 2 steps:
    1. Select the 2-best delays from N-best GCC-PHAT peaks in each individual channel.
    2. Select best delays vector across all channels.



# MDM-specific Improvements (II)

- Modified delay-sum algorithm (II)
  - Modified channel weighting ( $\alpha = 0.95$ )

$$W_i[n] = \alpha W_i[n-1] + (1 - \alpha) Xc_i[n]$$

Before:

$$Xc_i[n] = \frac{Xcorr(i, ref)[n]}{\sum_j Xcorr(j, ref)[n]}$$

Now:

$$Xc_i[n] = \frac{\sum_k Xcorr(i, k)[n]}{\sum_j \sum_k Xcorr(j, k)}$$

- “Suspicious” frames elimination if  $Xc_i[n] < \frac{1}{3(Nc-1)}$

# MDM-specific Improvements (III)

- Using delay&sum delays as features in combination with acoustic features.
  - Models built using a weighted combination of two separate feature streams, modeled by 2 separate GMM models.

$$\log p(x_{ac}[i], x_{del}[i] | \Theta_{tot}) = \alpha \cdot \log p(x_{ac}[i] | \Theta_{ac}) + (1 - \alpha) \cdot \log p(x_{del}[i] | \Theta_{dels})$$

- It is used in the Viterbi segmentation and in the BIC models comparison.

$$\Delta BIC_{tot} = \alpha \cdot \Delta BIC_{aco} + (1 - \alpha) \cdot \Delta BIC_{del}$$

# Eval Results

## Conference Room - Spkr

<i>Cond.</i>	<i>System ID</i>	<i>%DER</i>	<i>Description</i>
MDM	<b>p-wdels</b>	<b>35.77</b>	Primary system.
	c-newspnspdelay	35.77	Same as last year, w/ delays and new spnsp.
	c-wdelsfix	38.26	Same as primary, but init clusts = 16.
	c-nodels	41.93	Same as primary, but no delay feats.
	c-oldbase	42.36	Same as last year's system, w/ new spnsp.
SDM	<b>p-nodels</b>	<b>43.59</b>	Primary system (no delay feats).
	c-oldbase	43.93	Same as last year's system /w new spnsp.

Breakdown of errors by type for **primary systems**:

	<i>Miss</i>	<i>FA</i>	<i>SpNsp</i>	<i>Spkr</i>	<i>Total</i>
MDM	<b>27.60</b>	1.10	28.70	7.20	35.77
SDM	<b>28.90</b>	0.80	29.70	13.90	43.59

# Eval Results

## Lecture Room - Spkr

<i><b>Cond.</b></i>	<i><b>System ID</b></i>	<i><b>%DER</b></i>	<i><b>Description</b></i>
MDM	<b>p-wdels</b>	<b>24.01</b>	Primary system. (Same as conf. room system)
	c-nodels	23.63	Same as primary, but no delay feats.
	c-wdelsfix	24.53	Same as primary, but init # clusts = 10.
	c-guessone	26.96	Guess one speaker all the time. No sp/nsp.
SDM	<b>p-nodels</b>	<b>23.95</b>	Primary system (no delay feats).
	c-guessone	26.96	Guess one speaker all the time. No sp/nsp.
ADM	<b>p-wdels</b>	<b>21.05</b>	Same as MDM primary, but using all channels.
	c-nodels	20.24	Same as primary, but no delay feats.
	c-wdelsfix	22.11	Same as primary, but init # clusts = 10.
	c-guessone	26.96	Guess one speaker all the time. No sp/nsp.
MSLA	<b>p-guessone</b>	<b>26.96</b>	Guess one speaker all the time. No sp/nsp.

# Eval Results

## Conference Room - SAD

<b><i>Cond.</i></b>	<b><i>System ID</i></b>	<b><i>%Error</i></b>	<b><i>%Miss</i></b>	<b><i>%FA</i></b>	<b><i>Description</i></b>
MDM	p-dual	23.51	22.76	0.8	Two-pass system tuned to forced alignments.
SDM	p-dual	24.95	24.24	0.8	Two-pass system tuned to forced alignments.



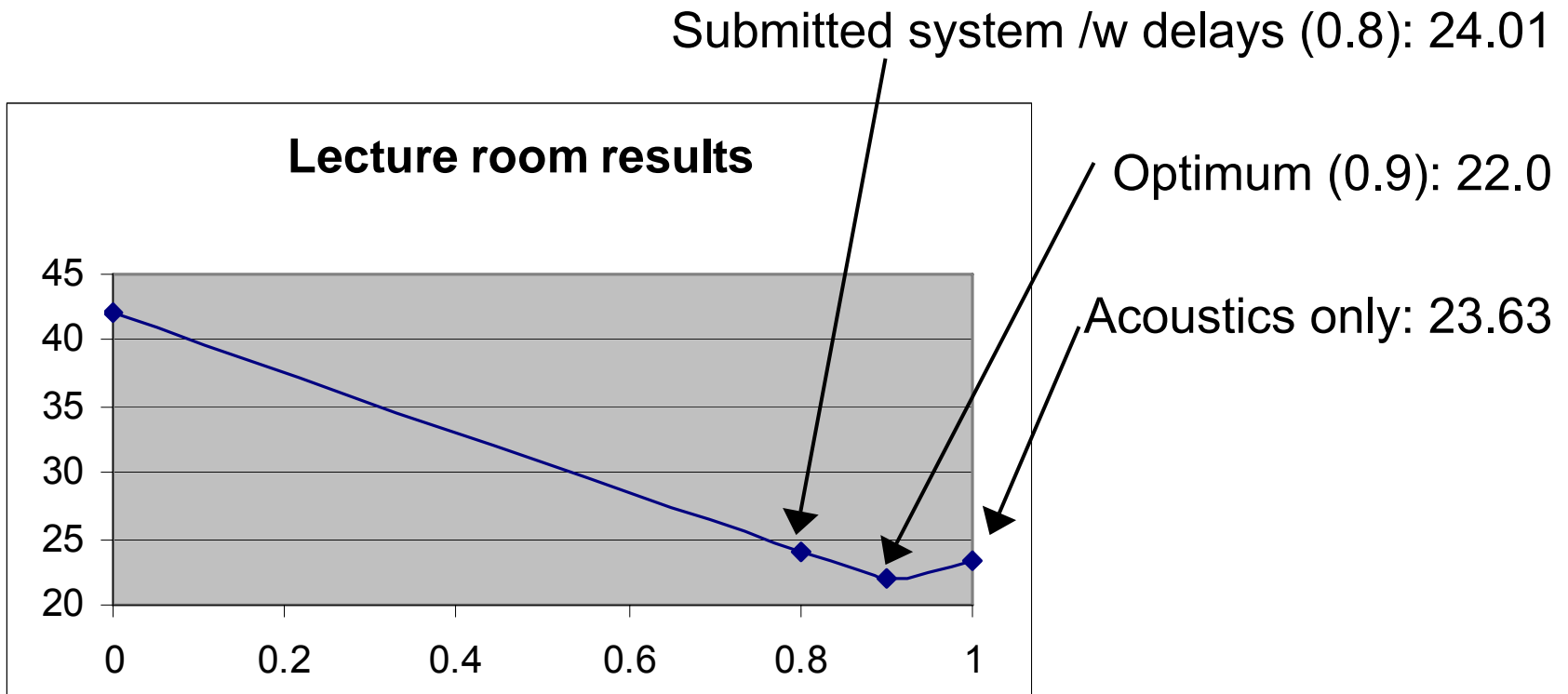
# Eval Results

## Lecture Room - SAD

<b><i>Cond.</i></b>	<b><i>System ID</i></b>	<b><i>%Error</i></b>	<b><i>%Miss</i></b>	<b><i>%FA</i></b>	<b><i>Description</i></b>
MDM	p-dual	13.83	9.3	4.5	Two-pass system tuned to forced alignments.
SDM	p-dual	14.59	10.0	4.6	Two-pass system tuned to forced alignments.
ADM	p-dual	13.22	9.3	3.9	Two-pass system tuned to forced alignments.

# Post-evaluation analysis

## ■ Lecture room



# Post-evaluation analysis

- Conference room data evaluated using ICSI-SRI Forced Alignments.

<i><b>Cond.</b></i>	<i><b>System ID</b></i>	<i><b>%DER FA</b></i>	<i><b>%DER hand</b></i>
MDM	<b>p-wdels</b>	<b>19.16</b>	<b>35.77</b>
	c-newspnspdelay	20.03	35.77
	c-wdelsfix	23.32	38.26
	c-nodels	27.46	41.93
	c-oldbase	27.01	42.36
SDM	<b>p-nodels</b>	<b>28.25</b>	<b>43.59</b>
	c-oldbase	28.21	43.93

# Future work

- Continue work on cluster initialization and purification.
- Add other acoustic features (e.g. PLP, prosody, etc.) as additional streams.
  - Try to improve both SDM and MDM/ADM conditions.
  - Dynamic stream weights.
- Experiment with alternatives to speech/non-speech, e.g. voiced/unvoiced.



# Questions?